

분류 및 예측 (1)

: 의사결정나무(Decision Tree)



분류/예측 기법 응용

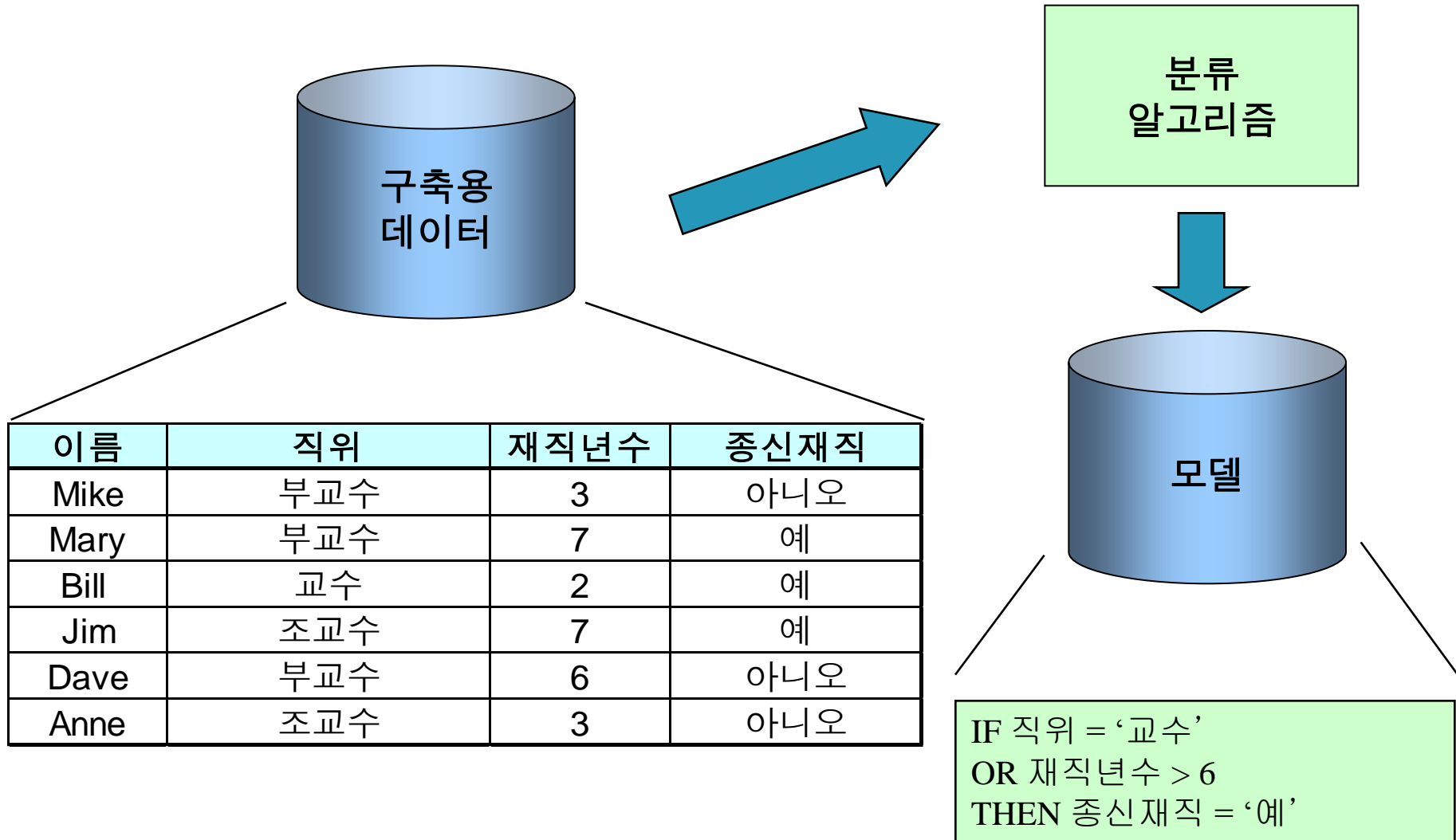
- 특정 상품(예: 화장품)의 구매 가능성이 높은 고객은 누구인가?
- 향후 6개월 안에 이탈할 가능성이 높은 고객들은 누구이며 그들의 특징은?
- 향후 1개월 안에 내점하여 최소 10만원 이상 구매할 가능성이 높은 고객은 누구인가?
- H백화점의 고객별 LTV를 계산해 주는 모델은?
 - 예: $LTV = \text{현재 고객의 기여가치} + \text{추가구매확률에 의한 기여가치} - \text{이탈확률에 의한 손실 가치}$
- 최고의 고객들은 어떤 특성을 갖는가?
- 최고 등급의 고객이 될 가능성이 높은 고객들은?



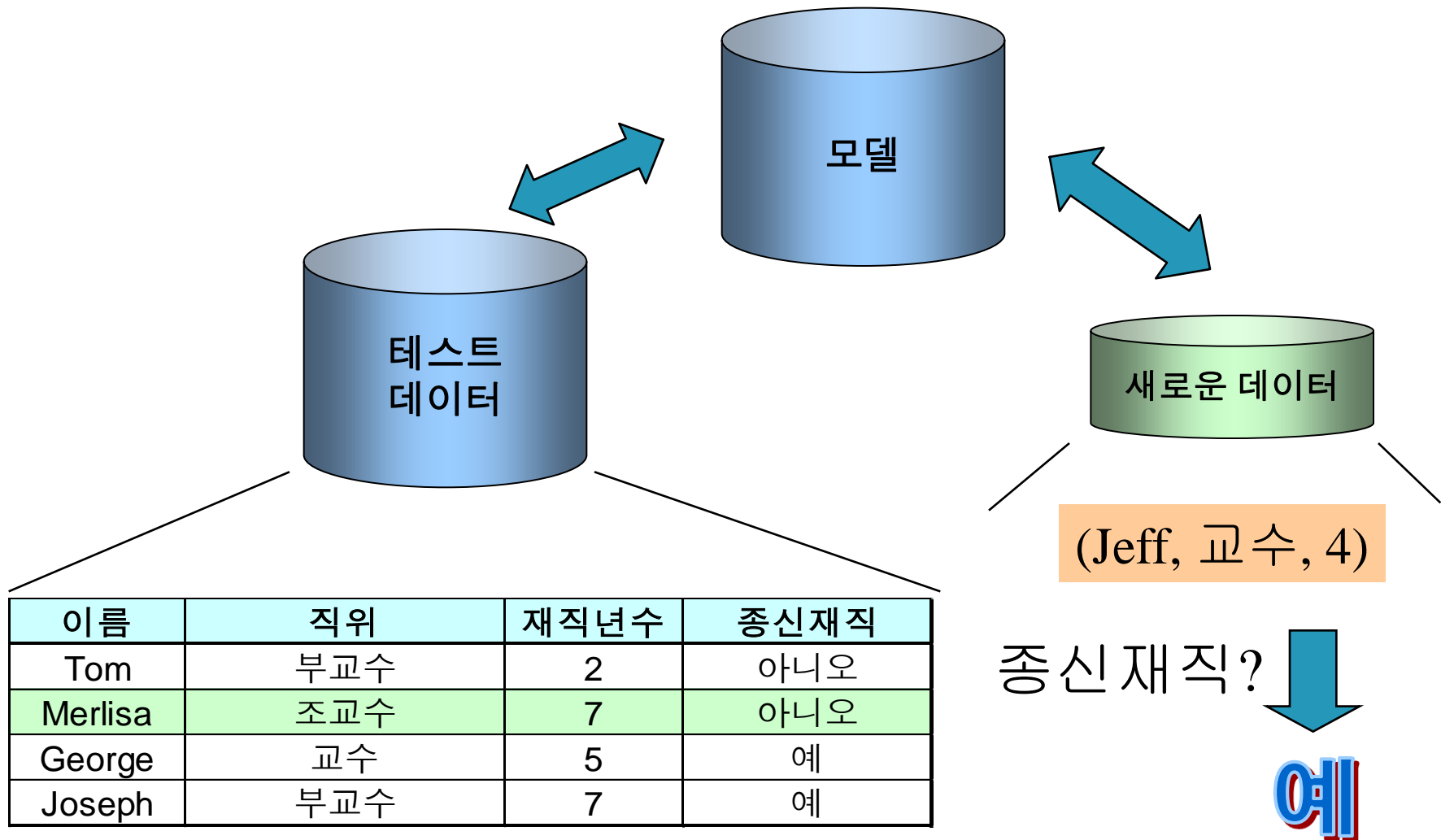
의사결정나무(Decision Tree)의 개요

- 데이터 분류와 예측에서 강력하고 인기 있는 도구
- 장점
 - 인간이 쉽게 이해할 수 있는 언어로 표현할 수 있는 규칙을 기반으로 판단
- 적용 예
 - 메일 마케팅 회사의 마케팅 모델
 - 회원들에 잘 반응하는 마케팅 모델을 예측하는데 사용
 - 은행에서 대출 상담의 경우
 - 대출 상담자에게 대출을 해 줄 경우, 정확한 근거를 제시하여 대출을 거부할 자료를 제시할 수 있음
 - If 연간 수입 \$20,000이상 and 관련 계좌가 3개 이상 then 대출
- 의사결정나무 기법
 - 명목형 목표변수 : C5.0, QUEST(Quick Unbiased Efficient Statistical Tree)
 - 연속형 목표변수 : CART(Classification & Regression Tree), CHAID(Chisquared Automatic Interaction Detection)

분류를 위한 모델 구축

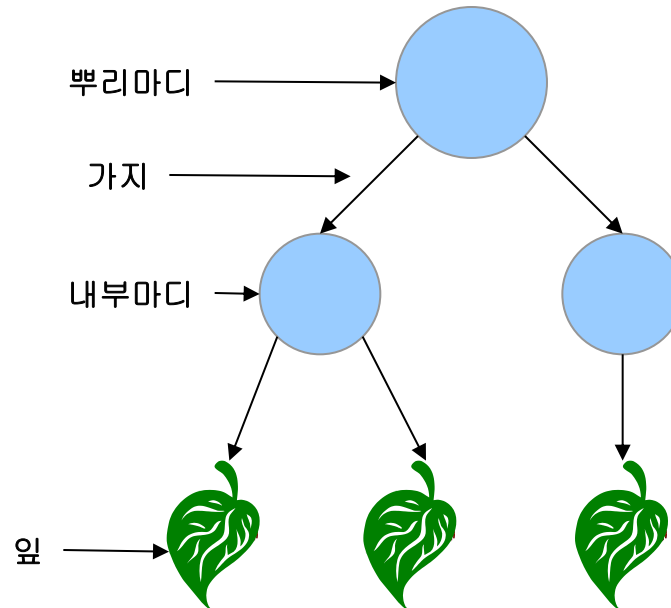


예측을 위한 모델의 이용



의사결정나무의 구성

- 뿌리마디(root node): 최상단에 위치
- 내부마디(internal node): 속성의 분리 기준을 포함
- 가지(link): 마디와 마디를 이어줌
- 잎(leaf): 최종 분류





의사결정나무의 형성과정

Step 1

의사결정나무 형성

- ✓ 분석의 목적과 자료구조에 따라, 적절한 최적의 분리기준(split criterion)을 찾아서 나무를 성장 시킨다. 정지규칙(stop rule) 을 만족하면 성장을 중단한다.

Step 2

가지치기

- ✓ 분류오류(classification error)를 크게 할 위험(risk)이 높거나 부적절한 추론규칙(induction rule)을 가지고 있는 가지(branch)를 제거한다. 또한, 불필요한 가지를 제거한다.

Step 3

타당성 평가

- ✓ 이익도표(gain chart)나 위험도표(risk chart) 또는 검증용 자료 (test sample)의 사용, 또는 교차타당성 (cross validation) 등을 이용하여 의사결정나무를 평가한다.

Step 4

해석 및 예측

- ✓ 구축된 나무모형을 해석하고 예측모형을 설정한다

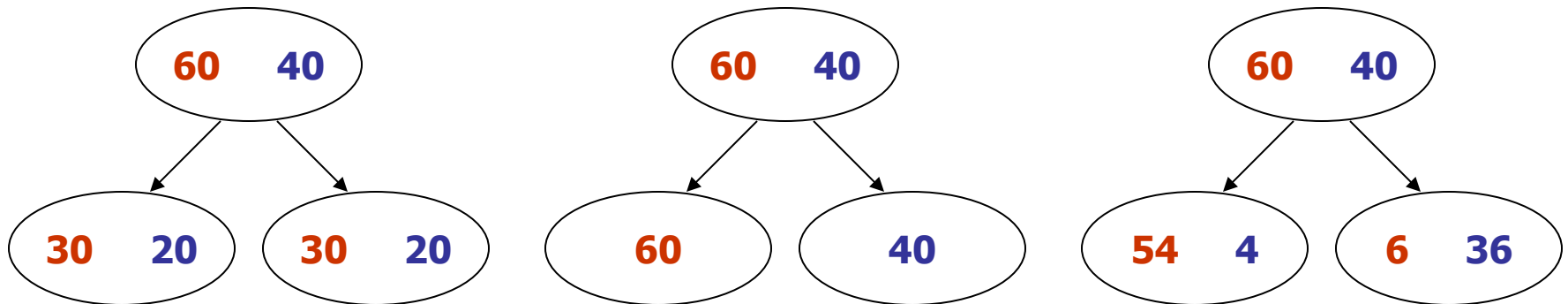
순수도 & 분리기준

■ 순수도

- 목표변수의 특정 범주에 개체들이 포함되는 정도

■ 분리기준

- 하나의 부모마디로부터 자식마디들이 형성될 때 입력변수의 선택과 범주의 병합이 이루어 질 기준을 의미
- 어떤 입력변수를 이용하여 분리하는 것이 목표변수의 분포를 가장 잘 구별해 주는지를 파악하여 자식마디가 생성되는데, 부모마디의 순수도에 비해서 자식마디들의 순수도가 증가하도록 자식마디를 형성





정지규칙 (stopping rule)

현재의 마디가 더 이상 분리가 일어나지 못하게 하는 규칙

- 모든 자료가 한 그룹에 속할 때
- 마디에 속하는 자료가 일정한 수 이하일 때
- 불순도의 감소량이 아주 적을 때
- 뿌리마디로부터 깊이가 일정 수 이상일 때



과도/과소 적합

- 과도적합(overfitting)

- 모델이 데이터에 필요이상으로 적합한 모델
- 데이터 내에 존재하는 규칙 뿐만 아니라 불완전한 레코드도 학습

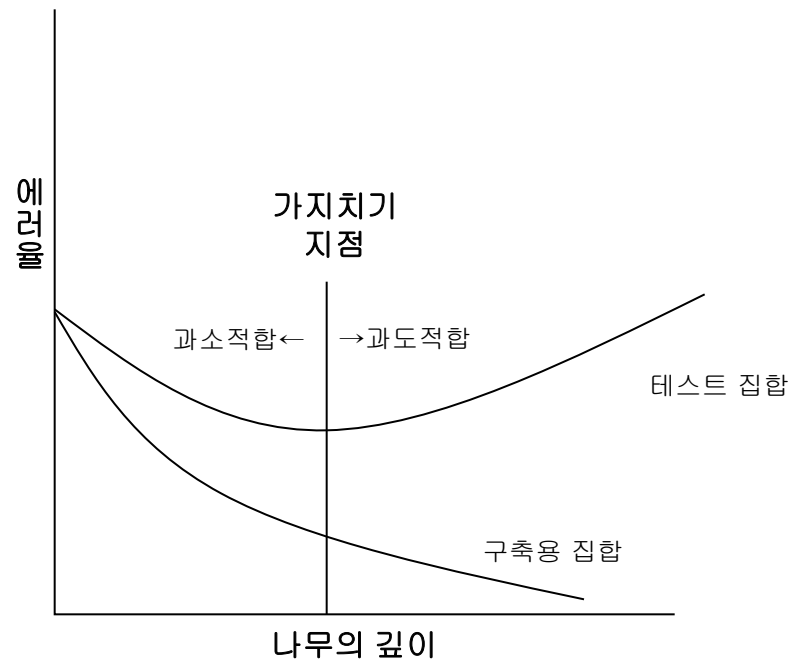
- 과소적합(underfitting)

- 모델이 데이터에 제대로 적합하지 못한 모델
- 데이터 내에 존재하는 규칙도 제대로 학습하지 못함

가지치기

■ 가지치기 규칙(pruning rule)

- 최종 마디의 수가 너무 많으면 모형이 과적합 상태가 됨 -> 현실문제에 적용할 수 있는 규칙이 나오지 않는다.
- 불필요하게 복잡해진 나무의 의미 없는 가지를 제거하는 작업



C5.0 알고리즘 (1/7)

C5.0은 킬란(Quilan, 1988)의 C4.5를 개량한 알고리즘으로 의사결정 나무 모델의 하나이며, 범주형(flag/set type)인 목표 변수를 이용하여 다른 입력변수들의 분류를 통해, 세분화 모델이나 목표변수를 예측할 수 있는 분석 방법이다.

C5.0의 규칙

✓ 예측변수의 선택기준

- ▶ 목표변수(종속변수)를 예측하는 데 있어 X_1, \dots, X_p 를 예측변수로 활용하고자 한다고 할때, C5.0처럼 나무구조의 결정규칙을 생성하기 위해서는 각 단계에서 p 개의 예측 변수 중 어느 것에 의하여 가지분리를 할 것인가를 선택해야 한다. 이 때 결정규칙들은 각기 다른 변수선택 기준을 쓴다.

✓ 정지규칙/가지치기(Pruning) 규칙

- ▶ 의사결정 나무 구조에서 사용되는 가지치기 규칙은 일반적으로 2 가지가 있다. 하나는 일정한 임계치를 기준으로 그보다 높은 경우에만 새로운 가지가 나오도록 하는 정지규칙이고, 다른 하나는 의사결정 나무가 다 자란 후 중요도가 비교적 작은 가지 분리들을 취소함으로써 나무의 크기를 감축시키는 가지치기규칙이다. C5.0은 가지치기 방법을 택하고 있으며, 이는 계산시간이 비교적 긴 반면 가지의 중요도를 확인한다는 점에서 큰 의미가 있다.

✓ 최소 레코드 수

- ▶ 가지 분리를 허용하는 조건으로 가지에 배속되는 레코드 수가 특정 값보다 커야 한다는 조건. 이 최소수에 대한 디폴트 값은 2인데, 이것을 큰 값으로 세팅할 수록 나무 규모가 줄어들게 된다. C5.0에서는 분리될 가지 중 2개 이상이 최소 레코드 수보다 커야 가지 분리가 허용된다.

✓ 엔트로피 지수 (Entropy Index) : C5.0에서 결정규칙을 분리할 수 있는 기준으로 데이터의 무질서 정도를 측정할 수 있는 방법

- ▶
$$Entropy(T) = - \sum_{i=1}^k p_i \log p_i \quad (p = \text{각 범주의 비율})$$

C5.0 알고리즘 (2/7)

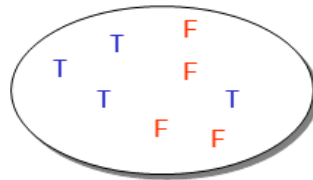
C5.0 -엔트로피(Entropy)

- ✓ C5.0 / C4.5의 기본 개념은 불순도를 측정하는(C&RT의 Gini와 유사)하나의 측정 지표가 엔트로피라는 것이다.
원래 엔트로피는 물리학의 용어로서 자연 현상의 변화는 일정한 방향으로만 진행되는, 즉, 자연현상은 물질계의 엔트로피가 증가하는 방향으로 진행된다는 엔트로피 증가의 법칙에서 나왔으며, 분자운동이 확률이 적은 질서 있는 상태에서부터 확률이 큰 무질서한 상태로 이동한다는 것을 의미한다.
- ✓ 이런 무질서와 확률을 이용한 개념을 데이터에서 구분을 잘 짓고, 못 짓는 속성 값을 찾아내는 지수로 사용하는 것으로 의사결정나무에서 활용을 한다.
- ✓ 의사결정 나무에서는 엔트로피가 높을 수록 Target 구분을 잘 못해주는 속성 필드가 되며, 낮을 수록 구분을 잘 해주는 유익한 속성 필드가 된다.

엔트로피 계산법

(Target 범주가 2개)

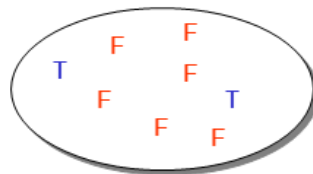
$$\text{Entropy}(S) = - \left(\frac{T \text{ 범주의 수}}{\text{전체건수}} \right) * \text{Log}_2 \left(\frac{T \text{ 범주의 수}}{\text{전체건수}} \right) - \left(\frac{F \text{ 범주의 수}}{\text{전체건수}} \right) * \text{Log}_2 \left(\frac{F \text{ 범주의 수}}{\text{전체건수}} \right)$$



-T범주수 : 4개, F범주수: 4개, 전체건수 8개인 경우

-실계산 : $-4/8 * \log_2(4/8) - 4/8 * \log_2(4/8) = 1$

※ 밑이 2인 로그임.



-T범주수 : 2개, F범주수: 6개, 전체건수 8개인 경우

-실계산 : $-2/8 * \log_2(2/8) - 6/8 * \log_2(6/8) = 0.5936$

※ 밑이 2인 로그임.

C5.0 알고리즘 (3/7)

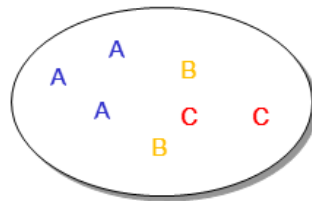
C5.0 – 엔트로피(Entropy)

- ✓ 엔트로피의 성질 : 집합에 범주가 고르게 분포할 수록 엔트로피 값은 높다. (불확실성 상태가 크다.)

엔트로피 계산법

(Target 범주가 3개)

$$\begin{aligned} \text{Entropy}(S) = & - (1\text{st 범주의 수} / \text{전체건수}) * \text{Log}(1\text{st 범주의 수} / \text{전체건수}) \\ & - (2\text{nd 범주의 수} / \text{전체건수}) * \text{Log}(2\text{nd 범주의 수} / \text{전체건수}) \\ & \dots \\ & - (N\text{th 범주의 수} / \text{전체건수}) * \text{Log}(N\text{th 범주의 수} / \text{전체건수}) \end{aligned}$$



A범주수: 3개, B범주수: 2개, C범주수: 2개, 전체건수 8개인 경우

$$\begin{aligned} \text{-실계산 : } & -3/8 * \log(4/8) - 2/8 * \log(2/8) - 2/8 * \log(2/8) \\ & = 0.4607 \end{aligned}$$

※ 밑이 10인 상용 로그임(주의 할 것).

C5.0 알고리즘 (4/7)

C5.0 – 정보획득함수(Information Gains)

- ✓ Information Gains라는 것은 C5.0에서 분류를 하는 지표가 되는 함수로서, CHAID의 Chi-square 통계량, C&RT, Gini Index와 같은 역할을 하게 된다.
- ✓ Information Gains라는 의미는 특정한 속성(attribute: 필드, 변수라고 생각하면 된다.)에 대한 정보를 알았을 때 얻어지는 정보량이 얼마만큼 되는지를 측정하는 지수가 되는 것이다.
- ✓ Entropy는 이 Information Gains를 계산하는데 필요한 측정식의 일부가 된다.

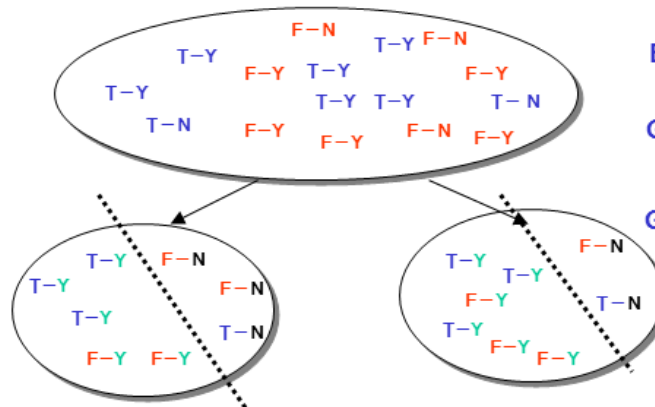
Information Gains

함수식

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

Information Gains

함수의 예제(개요)



$$Entropy(S) = 1$$

$$Gain(S, v1) = 1 - (3/8 * Entropy(N) + 5/8 * Entropy(Y)) = 0.048795$$

$$Gain(S, v2) = 1 - (2/8 * Entropy(N) + 6/8 * Entropy(Y)) = 0.271787$$



C5.0 알고리즘 (5/7)

C5.0 – 알고리즘 사례

ID	STR	INCOME	SEX	HOUSE	TARGET
1	서울	고소득	남자	아파트	구매
2	서울	고소득	남자	주택기타	구매
3	수도권	고소득	남자	아파트	비구매
4	지방	중간소득	남자	아파트	비구매
5	지방	저소득	여자	아파트	비구매
6	지방	저소득	여자	주택기타	구매
7	수도권	저소득	여자	주택기타	비구매
8	서울	중간소득	남자	아파트	구매
9	서울	저소득	여자	아파트	비구매
10	지방	중간소득	여자	아파트	비구매
11	서울	중간소득	여자	주택기타	비구매
12	수도권	중간소득	남자	주택기타	비구매
13	수도권	고소득	여자	아파트	비구매
14	지방	중간소득	남자	주택기타	구매

C5.0 알고리즘 (6/7)

C5.0 – 알고리즘 사례

1단계 : Target에 대한 Entropy를 계산한다.

→ Target 필드 (구매 5건, 비구매 9건) :

$$\text{Entropy} = -(5/14) * \log_2(5/14) - (9/14) * \log_2(9/14) = 0.9402$$

2단계 : Target을 제외한 다른 나머지 Input 필드(설명변수)의 Information Gains 값을 계산한다.

→ STR(지역)의 Information Gains 값 계산하기

1) "서울"범주의 Entropy(구매 3, 비구매 2) = $-(3/5)*\log_2(3/5) - (2/5)*\log_2(2/5) = 0.9709$

2) "수도권"범주의 Entropy(구매 0, 비구매 4) = $-(0/4)*\log_2(0/4) - (4/4)*\log_2(4/4) = 0$

3) "지방"범주의 Entropy(구매 2, 비구매 3) = $-(2/5)*\log_2(2/5) - (3/5)*\log_2(3/5) = 0.9709$

$$\text{Information Gains} = 0.9402 - (5/14)*0.9709 - (4/14)*0 - (5/14)*0.9709 = 0.2467$$

→ INCOME(소득)의 Information Gains 값 계산하기

1) "고소득"범주의 Entropy(구매 2, 비구매 2) = $-(2/4)*\log_2(2/4) - (2/4)*\log_2(2/4) = 1$

2) "중간소득"범주의 Entropy(구매 2, 비구매 4) = $-(2/6)*\log_2(2/6) - (4/6)*\log_2(4/6) = 0.9183$

3) "저소득"범주의 Entropy(구매 1, 비구매 3) = $-(1/4)*\log_2(1/4) - (3/4)*\log_2(3/4) = 0.8112$

$$\text{Information Gains} = 0.9402 - (4/14)*1 - (5/14)*0.9183 - (4/14)*0.8112 = 0.0948$$

→ SEX(성별)의 Information Gains 값 계산하기

1) "남자"범주의 Entropy(구매 4, 비구매 3) = $-(4/7)*\log_2(4/7) - (3/7)*\log_2(3/7) = 0.9852$

2) "여자"범주의 Entropy(구매 1, 비구매 6) = $-(1/7)*\log_2(1/7) - (6/7)*\log_2(6/7) = 0.5916$

$$\text{Information Gains} = 0.9402 - (7/14)*0.9852 - (7/14)*0.5916 = 0.1518$$

→ HOUSE(주거종류)의 Information Gains 값 계산하기

1) "아파트"범주의 Entropy(구매 2, 비구매 6) = $-(2/8)*\log_2(2/8) - (6/8)*\log_2(6/8) = 0.8112$

2) "주택기타"범주의 Entropy(구매 2, 비구매 4) = $-(3/6)*\log_2(3/6) - (3/6)*\log_2(3/6) = 1$

$$\text{Information Gains} = 0.9402 - (8/14)*0.8112 - (6/14)*1 = 0.048$$

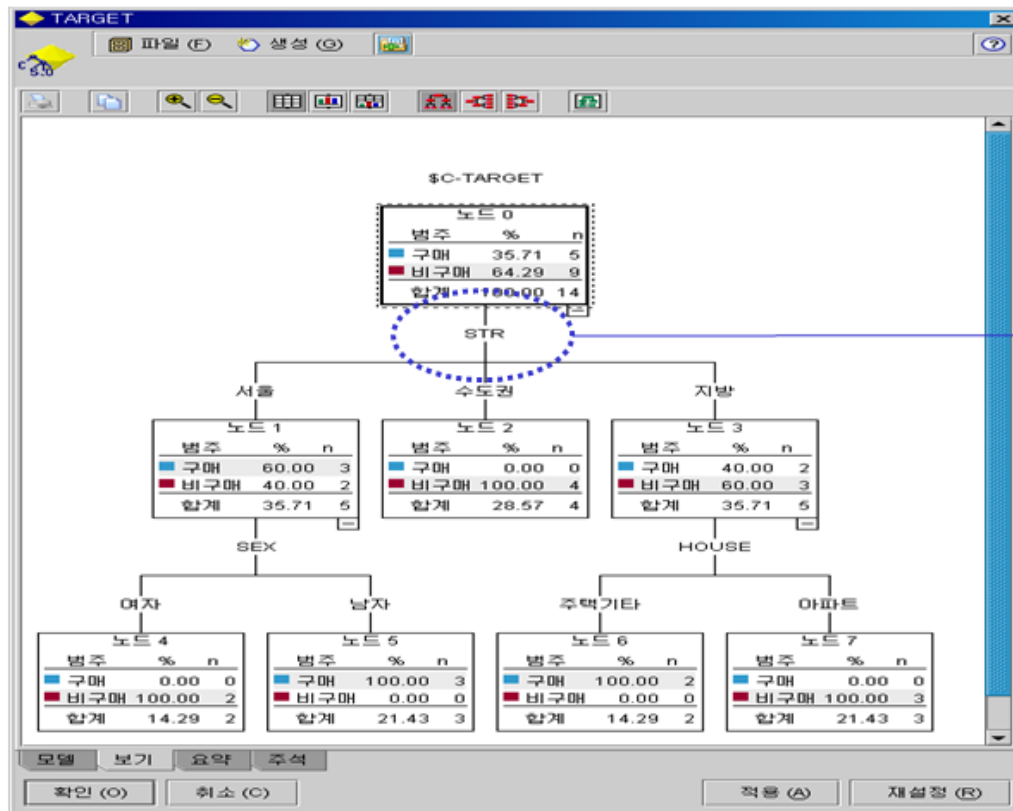
→ 가장 Information Gains 값이 큰 **STR**이 첫 번째 분류 기준 필드로 선정이 된다.

3단계 : 1단계 분류가 되면 자식노드(leaf)를 하나의 어미노드로 인식을 하여, 분류된 노드마다 2단계를 재 수행한다.

→ 단, STR 필드의 수도권 범주와 같이 0이 되면 더 이상 분류를 하지 않는다.(100% 만족된 상태)

C5.0 알고리즘 (7/7)

C5.0 - 알고리즘 사례



Clementine에서도 실질적으로 수행을 하면 STR 지역이 가장 먼저 분리가 되는 것을 확인할 수 있다.

C5.0 노드

C5.0

Fields Model Costs Analyze Annotations

Model name: ☐ Auto ☐ Custom

☒ Use partitioned data

☒ Build model for each split

Output type: ☒ Decision tree ☐ Rule set

☐ Group symbolics

☐ Use boosting Number of trials: 10

☐ Cross-validate Number of folds: 10

Mode: ☐ Simple ☒ Expert

① Pruning severity: 75

② ☒ Use global pruning

③ ☐ Winnow attributes

Minimum records per child branch: 2

OK Run Cancel Apply Reset

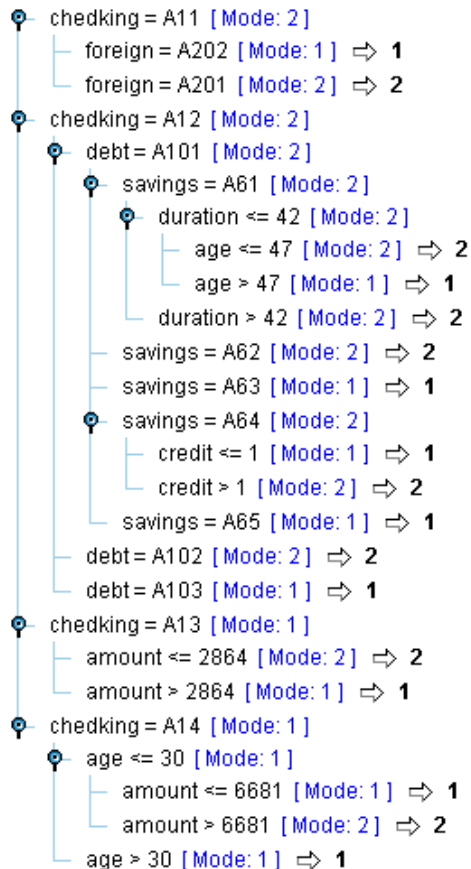
① C5.0의 문제점인 가지치기 강도를 조정하는 곳으로 이 값이 커지면 커질수록 가지치기 강도가 강해져서, Over-fitting의 가능성이 적어지지만, 대신 가지가 적게 되서, 정확도가 떨어지고, 이 값이 작을 수록 가지치기 가 약해 져서 Over-fitting할 확률이 높아지지만, 가지 가 많아져서 전체적인 정확도는 올라간다. default는 75%이다.

② Clementine C5.0에서는 2단계의 가지치기 (Pruning)를 하게 된다. 1단계는 지역(국소)적 가지 치기, 2단계는 전역적 가지치기가 그것인데 위의 가지 치기 강도는 지역적 가지치기의 강도 옵션 조정이라고 보면 된다. 이것은 해당되는 노드와 그 Sub-tree간의 가지치기가 이루어지는 것이고, 전역적 가지치기는 전체적으로 만들어진 Tree 구조에서 가지치기를 수행하는데 강도가 약한 sub-tree자체가 삭제되게 된다.

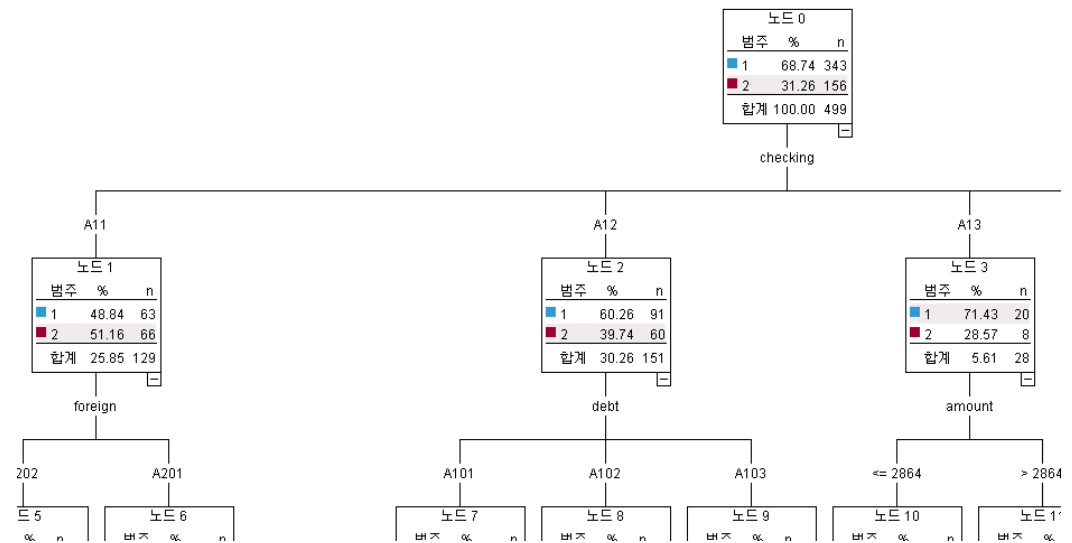
③ 입력 필드들에 대해서 사전에 필드가 유용한지 측정을 한 다음 유용하지 않는 경우 배제하고, 모델링을 수행한다.

C5.0 모형의 Output

Rule

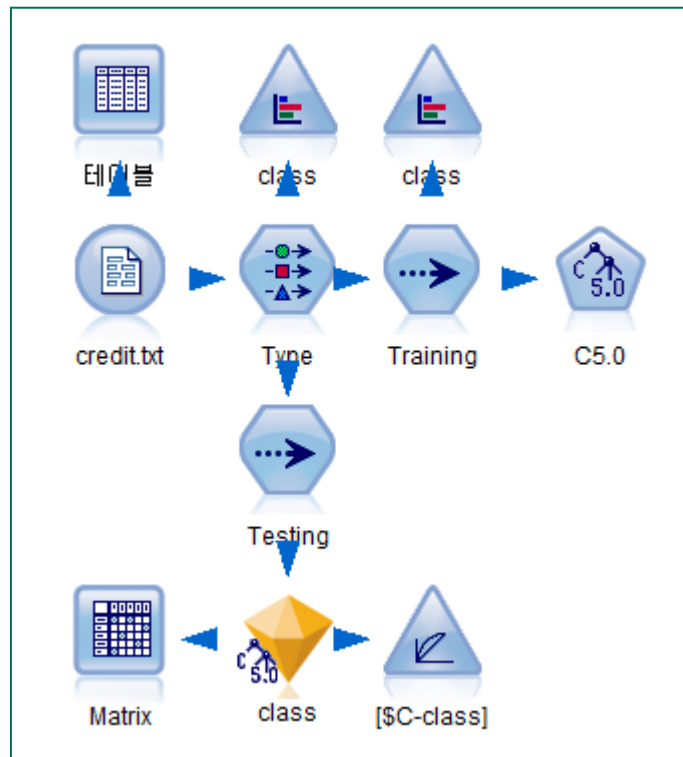


Tree



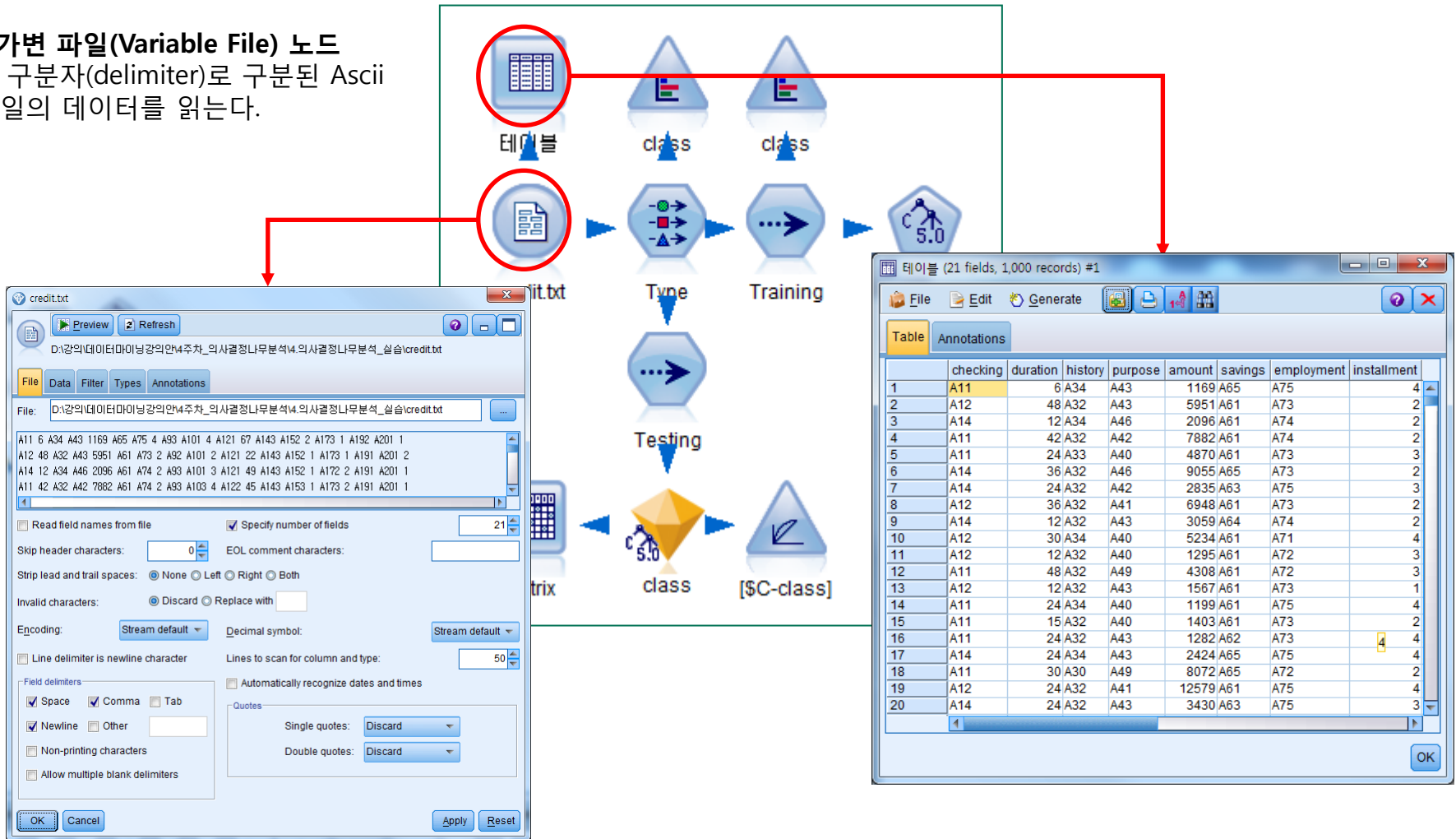
C5.0 실습 - 전체 스트림

신용 평가 데이터 credit.txt로 SPSS Modeler의 의사결정나무 분석을 실시



C5.0 실습 - (Step1) 데이터 연결 및 확인

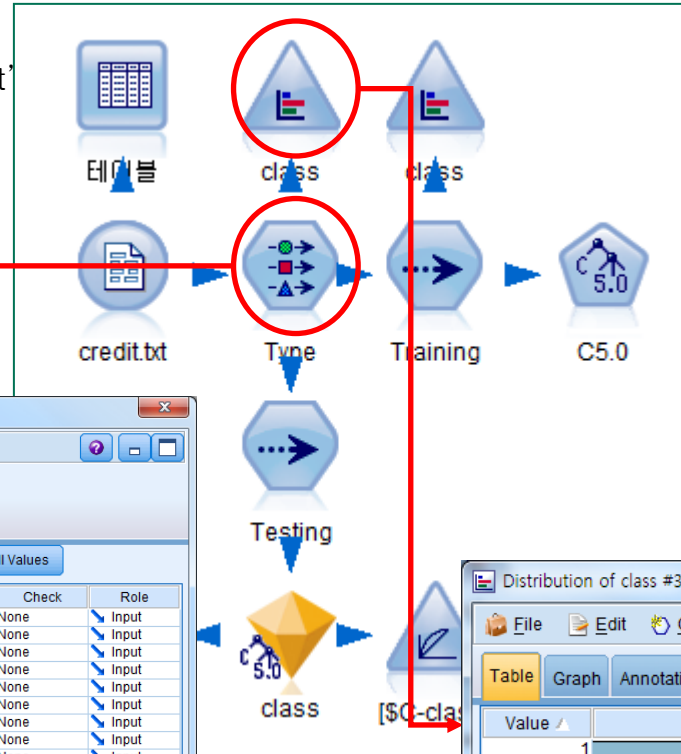
- 가변 파일(Variable File) 노드
: 구분자(delimiter)로 구분된 Ascii
파일의 데이터를 읽는다.



C5.0 실습 - (Step2) 데이터 탐색

▪ 유형 (Type) 노드

: Checking에서 Foreign까지는 'Input'으로 하고, Class는 'Target'으로 지정한다.



Type

Preview

Types Format Annotations

Read Values Clear Values Clear All Values

Field	Measurement	Values	Missing	Check	Role
checking	Nominal	A11,A12,A...	None	None	Input
duration	Continuous	[4,72]	None	None	Input
history	Nominal	A30,A31,A...	None	None	Input
purpose	Nominal	A40,A41,A...	None	None	Input
amount	Continuous	[250,18424]	None	None	Input
savings	Nominal	A61,A62,A...	None	None	Input
employment	Nominal	A71,A72,A...	None	None	Input
installment	Continuous	[1,4]	None	None	Input
marital	Nominal	A91,A92,A...	None	None	Input
debt	Nominal	A101,A102...	None	None	Input
residence	Continuous	[1,4]	None	None	Input
property	Nominal	A121,A122...	None	None	Input
age	Continuous	[19,75]	None	None	Input
plan	Nominal	A141,A142...	None	None	Input
housing	Nominal	A151,A152...	None	None	Input
credits	Continuous	[1,4]	None	None	Input
job	Nominal	A171,A172...	None	None	Input
people	Continuous	[1,2]	None	None	Input
telephone	Nominal	A191,A192...	None	None	Input
foreign	Nominal	A201,A202...	None	None	Input
class	Flag	2/1	None	None	Target

View current fields View unused field settings

OK Cancel Apply Reset

▪ 분포 (Distribution) 노드

: 종속변수 Class의 분포를 보기 위해 실행시키면 Class1(=good credit)이 700건(70%)이고, Class2(=bad credit)이 300건(30%)임을 알 수 있다.

Distribution of class #3

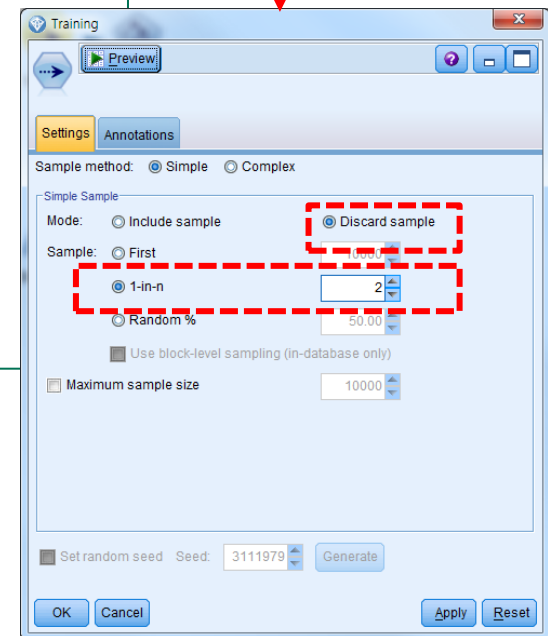
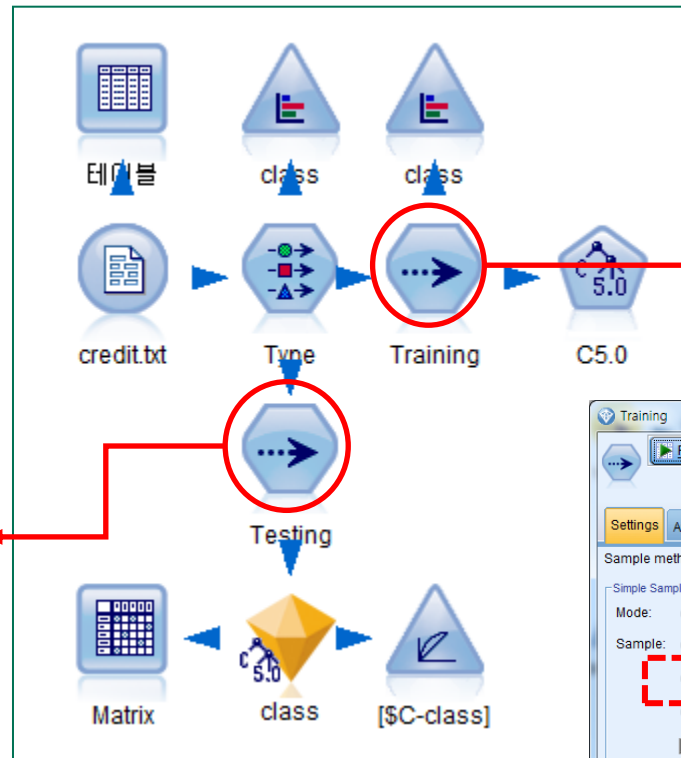
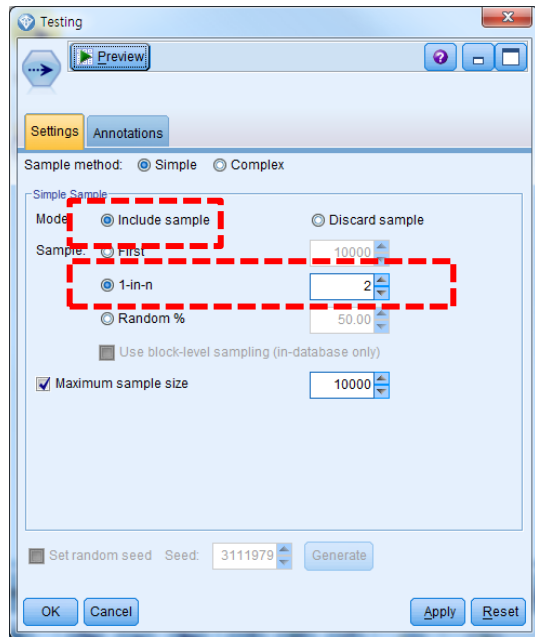
File Edit Generate View

Table Graph Annotations

Value	Proportion	%	Count
1		70.0	700
2		30.0	300

OK

C5.0 실습 - (Step3) 데이터 샘플링



C5.0 실습 - (Step4) 모델링

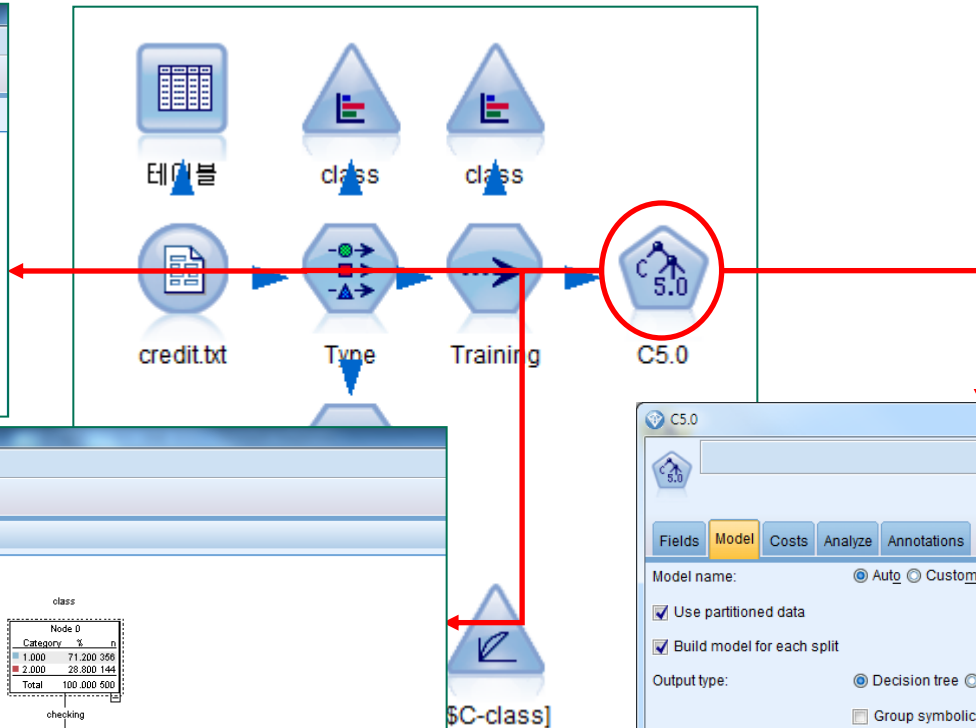
C5.0

File Generate View

Model Viewer Summary Annotations

```

checking = A11 [Mode: 2] (145)
├─ history in ["A30" "A31" "A32" "A33"] [Mode: 2] ⇒ 2 (112; 0.554)
├─ history in ["A34"] [Mode: 1] (33)
checking = A12 [Mode: 2] (118)
├─ property = A121 [Mode: 1] ⇒ 1 (38; 0.842)
├─ property = A122 [Mode: 2] (19)
├─ property = A123 [Mode: 2] (38)
├─ property = A124 [Mode: 2] ⇒ 2 (23; 0.652)
checking = A13 [Mode: 1] (35)
├─ residence <= 3 [Mode: 1] (26)
├─ residence > 3 [Mode: 1] ⇒ 1 (9; 1.0)
checking = A14 [Mode: 1] (202)
├─ plan = A141 [Mode: 2] (22)
├─ plan = A142 [Mode: 2] ⇒ 2 (5; 0.4)
├─ plan = A143 [Mode: 1] ⇒ 1 (175; 0.92)
    
```



C5.0

File Generate View

Model Viewer Summary Annotations

class

Node 0	Category	%	n
	1,000	71.200	356
	2,000	28.800	144
Total		100.000	500

checking

A12

Node 8	Category	%	n
	1,000	61.864	73
	2,000	38.136	45
Total		100.000	118

property

A13

Node 20	Category	%	n
	1,000	82.857	29
	2,000	17.143	6
Total		100.000	35

residence

A122

Node 10	Category	%	n
	1,000	52.632	23
	2,000	47.368	9
Total		100.000	32

residence

A123

Node 13	Category	%	n
	1,000	60.526	23
	2,000	39.474	15
Total		100.000	38

savings

A124

Node 19	Category	%	n
	1,000	34.783	8
	2,000	65.217	15
Total		100.000	23

residence

A121

Node 21	Category	%	n
	1,000	76.923	20
	2,000	23.077	6
Total		100.000	26

residence

A125

Node 26	Category	%	n
	1,000	100.000	9
	2,000	0.000	0
Total		100.000	9

residence

C5.0

Fields Model Costs Analyze Annotations

Model name: ☐ Auto ☐ Custom

☒ Use partitioned data

☒ Build model for each split

Output type: ☒ Decision tree ☐ Rule set

☐ Group symbolics

☐ Use boosting Number of trials:

☐ Cross-validate Number of folds:

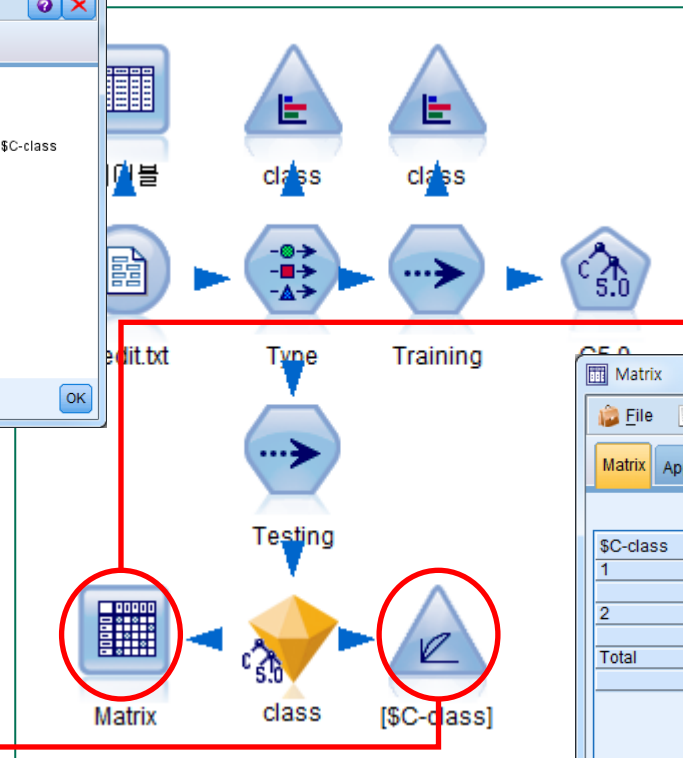
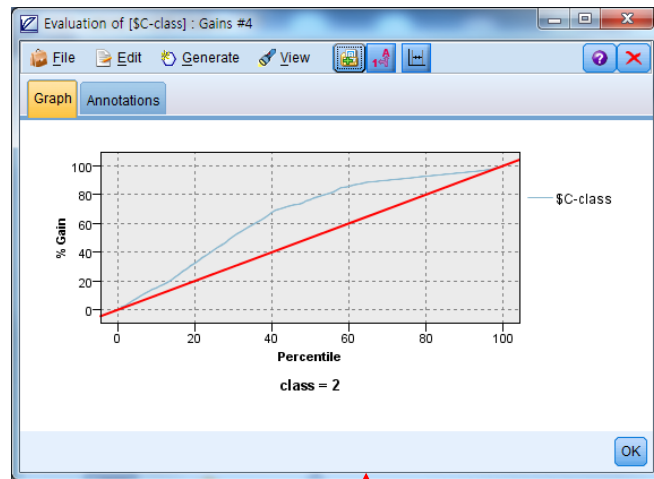
Mode: ☒ Simple ☐ Expert

Favor: ☐ Accuracy ☒ Generality

Expected noise (%):

OK Run Cancel Apply Reset

C5.0 실습 - (Step5) 모델 검증



Matrix

File Edit Generate

Matrix Appearance Annotations

class

\$C-class		1	2	Total
1	Count	221	37	258
1	Row %	85.659	14.341	100
2	Count	123	119	242
2	Row %	50.826	49.174	100
Total	Count	344	156	500
Total	Row %	68.800	31.200	100

Cells contain: cross-tabulation of fields (including missing values)

Chi-square = 70.581, df = 1, probability = 0

OK

모형 평가 – Gains Chart (1/2)

@ Gains

- 목표범주 1에 속하는 개체들이 각 등급에 얼마나 분포하고 있는지를 나타냄

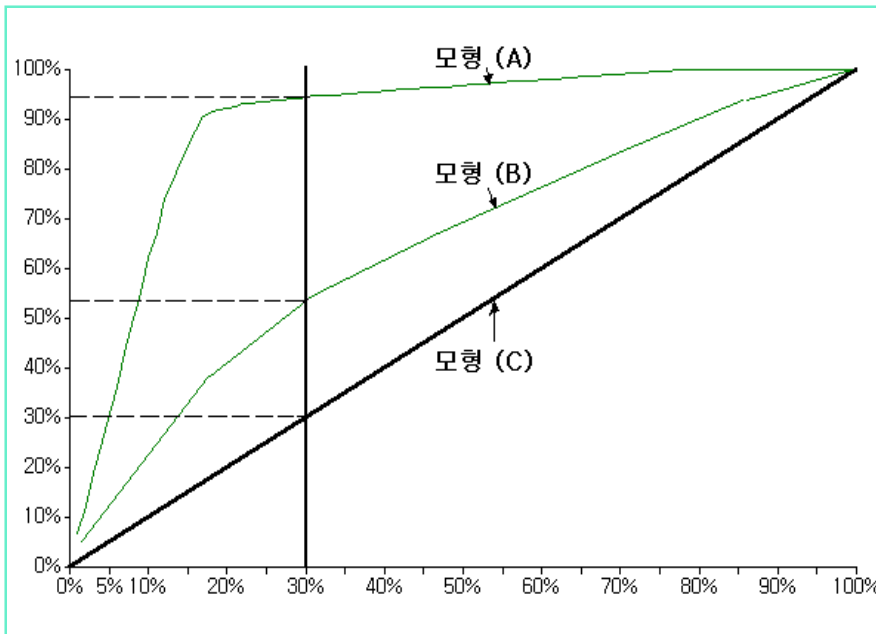
$$\frac{\text{해당 등급에서 목표변수의 특정 범주 빈도}}{\text{전체에서 목표변수의 특정 범주 빈도}} \times 100$$

등급	비누적				누적			
	빈도			반응률	빈도			반응률
	합계	Y=1	Y=0	Gain	합계	Y=1	Y=0	Gain
1	200	174	26	174/381=45.6	200	174	26	174/381=45.6
2	200	110	90	110/381=28.8	400	284	116	284/381=74.5
3	200	38	162	38/381= 9.9	600	322	278	322/381=84.5
4	200	14	186	14/381= 3.6	800	336	464	336/381=88.1
5	200	11	189	11/381= 2.8	1000	347	653	347/381=91.0
6	200	10	190	10/381= 2.6	1200	357	843	357/381=93.7
7	200	7	193	7/381= 1.8	1400	364	1036	364/381=95.5
8	200	10	190	10/381= 2.6	1600	374	1226	374/381=98.1
9	200	3	197	3/381= 0.7	1800	377	1423	377/381=98.9
10	200	4	196	4/381= 1.0	2000	381	1619	381/381=100.0
전체	전체 반응률 = 381/2000=19%							

모형 평가 - Gains Chart (2/2)

㉠ Gains Chart

- ㉠ 해당 등급에 따라 계산된 Gain값을 연속적으로 연결한 도표
- ㉠ 차트에서 볼 수 있는 좌하에서 우상을 걸친 대각선은 모형비교의 기준선으로서, 모형성능이 나쁘면 나뉠수록 이 기준선에 가까워짐



- ㉠ DM 발송에 대한 반응여부라는 목표변수의 '반응했음'이라는 범주에 대한 차트라고 가정해 보자
- ㉠ 서로 다른 방법으로 구축된 3개의 모형으로 동일한 수의 DM을 발송하는 경우, 얻어지는 반응률의 차이를 알 수 있음
- ㉠ 전체 관찰치 중 30%를 대상으로 DM을 발송하였을 때, 모형 (A)는 90%이상의 반응을 보임
- ㉠ 반면 모형 (B)는 50% 조금 넘는 반응을 보임
- ㉠ 따라서 분석자는 모형 (A)를 선택하게 됨

Case Study – 반품고객 예측 (1/4)

상황

국내 홈쇼핑 A사는 최근 소비자의 반품 횟수가 증가됨에 따라 마케팅 부서의 김팀장이 반품고객의 특성을 파악하고자 함.

데이터

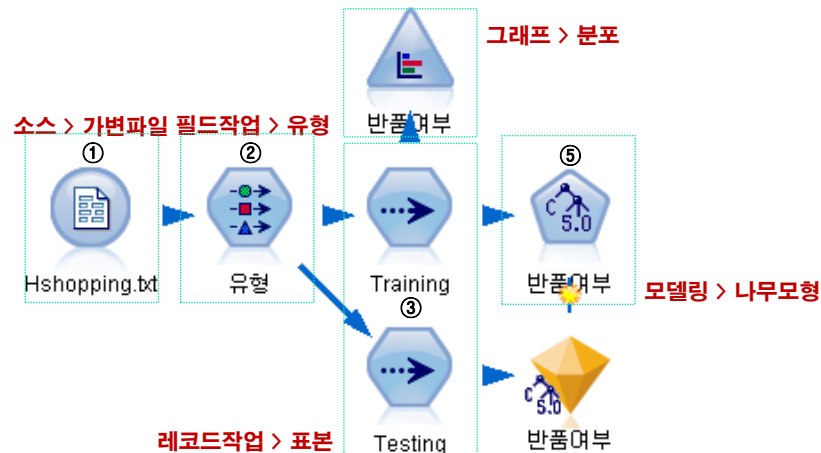
홈쇼핑 A사 고객 500명에 대한 성별, 나이, 구매 금액, 홈쇼핑 출연자, 반품 여부

분석 과정

① 데이터 준비 → ② 변수 지정 → ③ 훈련 · 테스트자료 분류 → ④ 균형화 작업 → ⑤ 의사결정 나무분석

Data: Hshopping.txt

의사결정 나무분석 과정



Case Study – 반품고객 예측 (2/4)

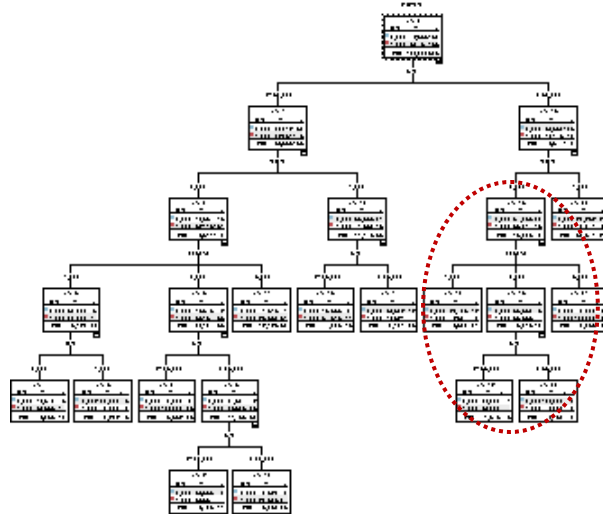
▷ 변수 목록

No.	변수 이름		변수 설명	변수 유형
	SPSS용	SAS용		
1	ID	ID	고객 고유번호	수치형
2	성별	SEX	1=남자, 2=여자	범주형
3	나이	AGE	나이	수치형
4	구매금액	BUYM	1=10만 원 미만, 2=10~30만 원, 3=30만 원 이상	범주형
5	출연자	ACTOR	1=일반인, 2=유명인	범주형
6	반품 여부	RETURNSYN	0=반품 ×, 1=반품 ○	범주형

Case Study – 반품고객 예측 (3/4)

▷ 실습 결과

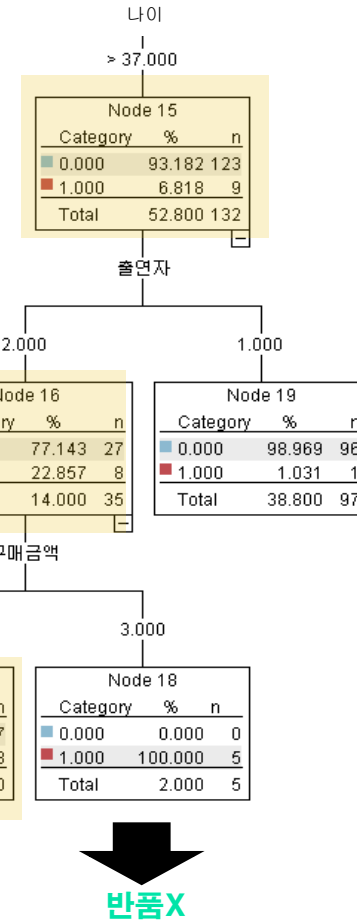
▪ 의사결정나무 분석 모형 > 뷰어



부분
확대

▪ 규칙 (부분)

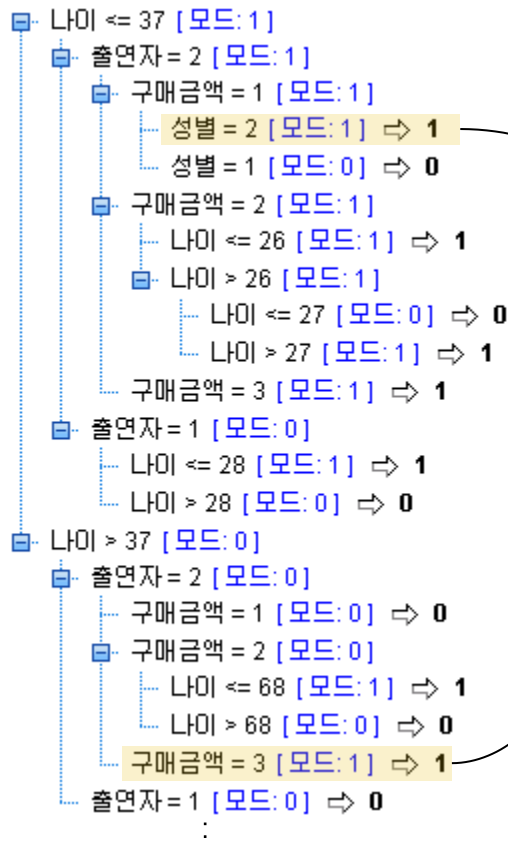
출연자가 유명인 & 구매금액이 10만원~30만원 &
나이가 68세 이하 → 반품O



Case Study – 반품고객 예측 (4/4)

▷ 실습 결과

▪ 의사결정나무 분석 모형 > 모델



▪ 규칙 1

나이가 37세 이하 & 출연하는 사람이 유명인 & 구매 금액 10만원 미만 & 성별이 여성

→ 반품 O

▪ 규칙 2

나이가 37세 초과 & 출연하는 사람이 유명인 & 구매 금액 30만원 이상

→ 반품 O